



4<sup>th</sup>

第四届博鳌国际药械真实世界研究大会

BOAO INTERNATIONAL CONFERENCE ON REAL WORLD STUDIES OF MEDICAL PRODUCTS

## 大型前列腺癌临床数据集向 OMOP 通用数据模型的转换 – 基于中国区域医疗数据库的经验

Transforming a Large-Scale Prostate Cancer Outcomes Dataset  
to the OMOP Common Data Model

- Experiences with a Regional Health Record Database in China

汇报人：许哲 单位：拜耳医药保健有限公司

Presenter: Zhe Xu Affiliation: Bayer Healthcare Company Ltd

代表项目作者汇报 On behalf of co-authors





## 4<sup>th</sup> 第四届博鳌国际药械真实世界研究大会

BOAO INTERNATIONAL CONFERENCE ON REAL WORLD STUDIES OF MEDICAL PRODUCTS

**作者：蔡启亮<sup>1</sup>, 谢洋<sup>2</sup>, 李静<sup>2</sup>, 林晓羽<sup>2</sup>, 杨芷其<sup>2</sup>, 殷铮<sup>2</sup>, 叶怡君<sup>2</sup>, 黄诗韵<sup>3</sup>, 许哲<sup>3</sup>, 牛远杰<sup>4</sup>**

<sup>1</sup>天津医科大学第二医院泌尿外科, 天津, 中国

<sup>2</sup>艾昆纬企业管理咨询有限公司真实世界研究, 上海, 中国

<sup>3</sup>拜耳医药保健有限公司医学事务部, 北京, 中国

<sup>4</sup>天津医科大学总医院泌尿外科, 天津, 中国

**Authors: Qiliang Cai<sup>1</sup>, Yang Xie<sup>2</sup>, Jing Li<sup>2</sup>, Xiaoyu Lin<sup>2</sup>, Chih-Chi Yang<sup>2</sup>, Zheng Yin<sup>2</sup>, Yi-Chun Yeh<sup>2</sup>, Shiyun Huang<sup>3</sup>, Zhe Xu<sup>3</sup>, Yuanjie Niu<sup>4</sup>**

<sup>1</sup>Department of Urology, Second Hospital of Tianjin Medical University, Tianjin, China.

<sup>2</sup>Real World Solutions, IQVIA Solutions Enterprise Management Consulting (Shanghai) Co., Ltd, Shanghai, China

<sup>3</sup>Medical Affairs, Pharmaceuticals, Bayer Healthcare Company Ltd, Beijing, China.

<sup>4</sup>Department of Urology, Tianjin Medical University General Hospital, Tianjin, China.

### • 讲者简介 Speaker Bio

许哲, 博士, 现任拜耳医药保健有限公司医学部证据生成高级经理, 负责多个产品的证据生成计划制定, 专病数据库搭建及数据标准化项目, 真实世界研究等。

先后于北京大学、美国约翰霍普金斯大学、英国剑桥大学获得预防医学学士、流行病学硕士、博士学位, 主要研究方向为基于大规模电子病历数据的心血管疾病风险预测模型。

**Zhe Xu, PhD.** I currently work as the Senior Manager of Evidence Generation at Bayer HealthCare Co. Ltd., mainly responsible for the evidence generation plans, establishment of real-world data cohorts, database standardization, and real-world studies, for multiple products.

I have obtained a Bachelor's degree in Preventive Medicine from Peking University, a Master's degree in Epidemiology from Johns Hopkins University, and a Ph.D. in Epidemiology from the University of Cambridge. My main research focus is on cardiovascular disease risk prediction models based on large-scale electronic health record data.



许哲 Zhe Xu

拜耳医药保健有限公司  
Bayer Healthcare Company Ltd.



## 由OHDSI维护的OMOP-CDM为真实世界数据分析提供一个标准化框架

OMOP-CDM, maintained by OHDSI, provides a standardized framework for real-world data analysis



**4<sup>th</sup> BOAO** INTERNATIONAL CONFERENCE ON  
**REAL WORLD STUDIES**  
OF MEDICAL PRODUCTS

第四届博鳌国际药械真实世界研究大会



### 实现快速且可靠的研究 Enables rapid and reliable research

标准化使研究人员能够在各种数据集和数据类型中快速而可靠地开展研究

Standardization allows researchers to conduct reliable studies across a variety of datasets and data types.



### 降低维护成本 Reduces the cost of ownership

通过减少理解数据模式和编写跨数据库分析代码所需的工作量来降低成本

Standardization lowers costs by reducing the effort required to understand data schemas and to write cross-database analytical code.



### 扩展数据访问 Expands data accessibility

OMOP-CDM 借助 OHDSI 网络的分布式多中心研究实现更广泛的数据访问

The OMOP-CDM facilitates expanded data access via the OHDSI network and remote multi-site database studies.



### OHDSI collaborators



Abbreviation: OHDSI = Observational Health Data Sciences and Informatics; OMOP-CDM = Observational Medical Outcomes Partnership Common Data Model

Reference:

OHDSI Collaborators. Accessed from <https://www.ohdsi.org/who-we-are/collaborators/>



## OMOP 被全球各国政府机构采用

OMOP is adopted by government agencies globally



**4<sup>th</sup> BOAO** INTERNATIONAL CONFERENCE ON  
**REAL WORLD STUDIES**  
OF MEDICAL PRODUCTS

第四届博鳌国际药械真实世界研究大会





## 国家药品监督管理局 (NMPA) 的指南中引用了通用数据模型 (CDM)

The Common Data Model (CDM) is cited in the guideline from the National Medical Products Administration (NMPA)

- **《用于产生真实世界证据的真实世界数据指导原则（试行）》**
  - 2021年4月15日发布
  - 由国家药品监督管理局药品审评中心 (CDE) 发布
  - 自发布之日起施行
- **Guidance for Real-World Data Used to Generate Real-World Evidence (Trial Implementation)**
  - Date of Issuance: April 15, 2021
  - Issuing Authority: Center for Drug Evaluation (CDE) of the National Medical Products Administration (NMPA)
  - Effective date: Effective immediately upon publication (April 15, 2021)



The screenshot shows the official website of the Center for Drug Evaluation (CDE) under the National Medical Products Administration (NMPA). The header includes the NMPA CDE logo and name in Chinese and English. Below the header, the current location is indicated as '新闻中心 >> 工作动态 >> 通知公告 >> 新闻正文'. The main content area features a title: '国家药监局药审中心关于发布《用于产生真实世界证据的真实世界数据指导原则（试行）》的通告（2021年第27号）' (Notice of the Center for Drug Evaluation regarding the issuance of the 'Guidelines for Real-World Data Used to Generate Real-World Evidence (Trial Implementation)' (2021年第27号)). The issuance date is listed as '发布日期: 20210415'. The body text explains the purpose of the guidelines and references the 'Guidelines for the Issuance of Technical Guidelines for Drug Evaluation' (药监综药管〔2020〕9号). The notice is signed by the '国家药品监督管理局药品审评中心' (Center for Drug Evaluation) on '2021年4月13日' (April 13, 2021). At the bottom, there is a section for attachments, showing '附件 1: 《用于产生真实世界证据的真实世界数据指导原则（试行）》.pdf'.

Abbreviation: NMPA = National Medical Products Administration; CDE = Center for Drug Evaluation

Reference:

1. Guidance for Real-World Data Used to Generate Real-World Evidence (Trial Implementation)

Accessed from <https://www.cde.org.cn/main/news/viewInfoCommon/2a1c437ed54e7b838a7e86f4ac21c539>.





## 国家药品监督管理局 (NMPA) 的指南中引用了通用数据模型 (CDM)

The Common Data Model (CDM) is cited in the guideline from the National Medical Products Administration (NMPA)



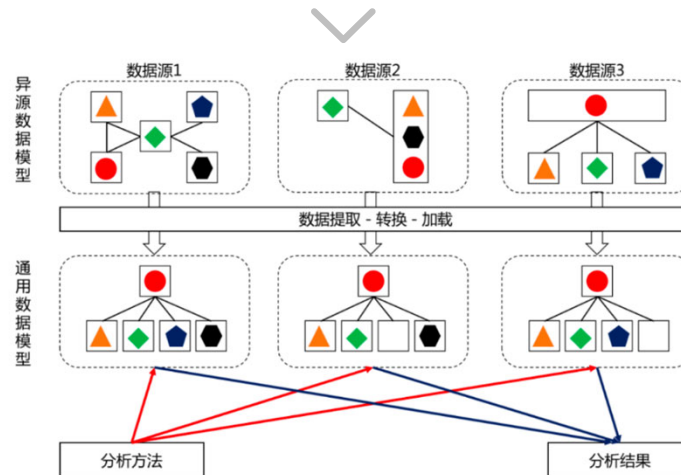
4<sup>th</sup> BOAO INTERNATIONAL CONFERENCE ON  
REAL WORLD STUDIES  
OF MEDICAL PRODUCTS

第四届博鳌国际药械真实世界研究大会

由于多源数据的结构和类型的复杂性、样本规模和标准的差异性，在将源数据转换为通用数据模型的整体过程中，需要对源数据进行提取、转换、加载，应确保源数据在语法和语义上与目标分析数据库的结构和术语一致。<sup>1</sup>

Given the complexity of multi-source data structures and types, and the heterogeneity in sample sizes and standards, converting source data into a common data model requires extraction, transformation, and loading. Source data must be syntactically and semantically aligned with the structure and terminology of the target analytical database<sup>1</sup>.

四、真实世界数据治理 .....	15
(一) 个人信息保护和数据安全性处理 .....	16
(二) 数据提取 .....	16
(三) 数据清洗 .....	17
(四) 数据转化 .....	18
(五) 数据传输和存储 .....	18
(六) 数据质量控制 .....	18
(七) 通用数据模型 .....	19
(八) 真实世界数据治理计划书 .....	21



Abbreviation: NMPA = National Medical Products Administration; CDE = Center for Drug Evaluation

Reference:

1. Guidance for Real-World Data Used to Generate Real-World Evidence (Trial Implementation)

Accessed from <https://www.cde.org.cn/main/news/viewInfoCommon/2a1c437ed54e7b838a7e86f4ac21c539>.

## 数据来源

- 截至2024年，天津电子健康档案数据库 (rEHR) 包含来自公立医疗机构的超过1100亿条记录，其中290亿条记录已通过数据治理处理，覆盖自1995年以来超过1600万居民。
- 数据整合自**43家三级医院**和**39家二级医院**，以及公共卫生记录，覆盖自2000年以来约1200万人。
- 项目提取了**2017年至2021年的数据**，并转换为 OMOP-CDM 模型，研究级数据质量自2015年起可用。

## Data source

- As of 2024, the Tianjin regional electronic health record database (rEHR) includes over 110 billion records from public medical institutions, with 29 billion records processed through data governance, covering over 16 million residents since 1995.
- It integrates data from **43 tertiary hospitals** and **39 secondary hospitals**, along with public health records, covering ~12 million people since 2000.
- For this study, **data from 2017 to 2021** were extracted and transformed into the OMOP-CDM schema, with research-grade data quality available since 2015.

## 队列定义 Cohort definitions

### 样本选择

- 纳入的患者在**天津电子健康档案数据库 (rEHR)** 的门诊或住院记录中**至少有一次前列腺癌 (PC) 诊断**。
- PC诊断通过ICD-10代码C61或诊断名称包含“前列腺”，并伴随“癌”、“原位癌”、“CA”、“恶性肿瘤”、“PCA”、“CRPC”或“HSPC”等术语进行识别。
- 如诊断名称包含“OPCA”、“PCAD”、“前列腺增生”或“肾癌”等术语，则排除，以确保病例识别的准确性。

### Sample selection

- Included **patients had at least one diagnosis of Prostate Cancer (PC)** recorded during outpatient visits or hospitalizations in the **Tianjin rEHR database**.
- PC diagnoses were identified using ICD-10 code C61 or diagnosis names containing “前列腺” along with terms like “癌”，“原位癌”，“CA”，“恶性肿瘤”，“PCA”，“CRPC”，or “HSPC”.
- Records were excluded if diagnosis names included terms such as “OPCA”，“PCAD”，“前列腺增生”，or “肾癌” to ensure accurate case identification.





## OMOP-CDM 转换工作流程

### OMOP-CDM conversion workflow



**4<sup>th</sup> BOAO** INTERNATIONAL CONFERENCE ON  
**REAL WORLD STUDIES**  
OF MEDICAL PRODUCTS  
第四届博鳌国际药械真实世界研究大会

#### 源数据探查 Source Data Profiling

1

- 分析源数据的结构和内容，以了解其格式、完整性和质量。
- Analyze the structure and content of the source data to understand its format, completeness, and quality.

#### ETL 映射到 CDM模式 ETL Mapping to CDM

2

- 使用提取、转换和加载 (ETL) 操作将源数据元素映射到 CDM 模式的指定目标位置。
- Map source data elements to the CDM schema using extract, transform, and load(ETL) operations.

#### 术语标准化 Vocabulary Standardization

3

- 将源数据中使用的代码和值转换为 CDM 中采用的代码和值，以实现数据的一致性与互操作性。
- Translate source codes and values into standardized CDM vocabularies for consistency and interoperability.

#### 质量保证 Quality Assurance

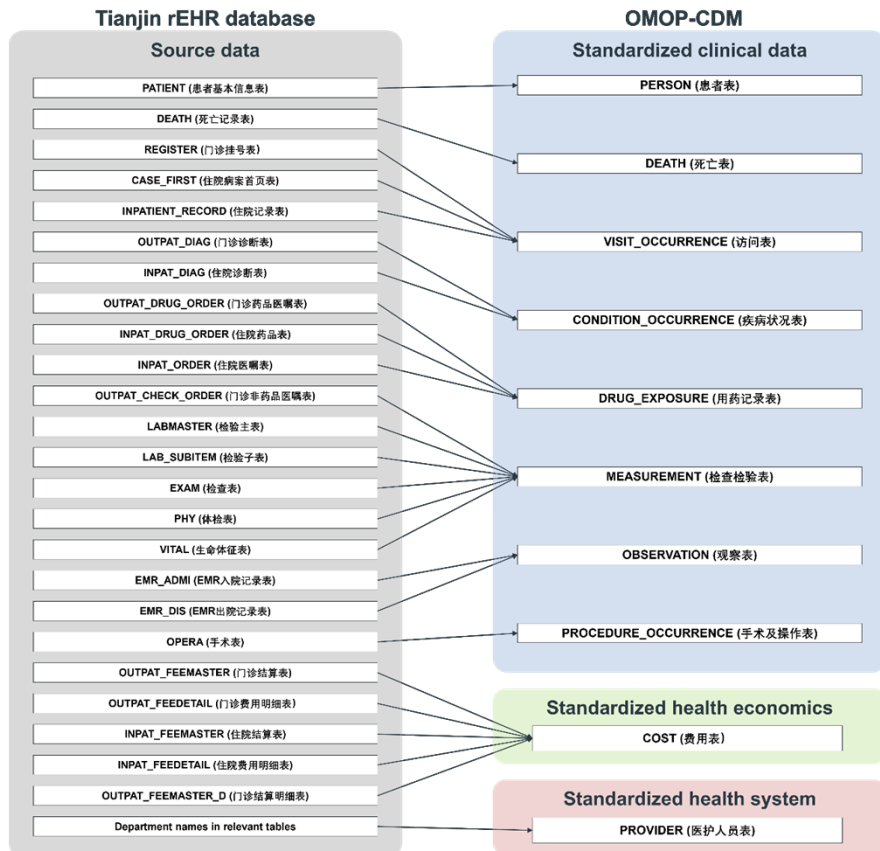
4

- 执行验证和质量检查，以确保数据的准确性、完整性以及符合 CDM 标准。
- Perform validation and quality checks to ensure data accuracy, completeness, and conformance to CDM standards.



## 将天津电子健康档案数据库转换为OMOP-CDM

### Transformation of the Tianjin rEHR Database into OMOP-CDM



天津电子健康档案数据已成功映射到 OMOP-CDM 标准化表格，确保后续分析和研究的互操作性和一致性。

The Tianjin rEHR data have been successfully mapped to OMOP-CDM standardized tables, ensuring interoperability and consistency for downstream analytics and research.



天津电子健康档案数据库 (rEHR) 提供 OMOP-CDM 所需的**所有**临床信息，涵盖**三个**主要类别：  
The Tianjin rEHR database provides **all** clinical information required for OMOP-CDM, covering **three** major categories:

- 临床数据 Clinical Data
- 健康经济学 Health Economics
- 卫生系统数据 Health System Data

# 天津电子健康档案数据库成功完成OMOP-CDM映射和质量控制 Successful OMOP-CDM Mapping and Quality Assurance in Tianjin rEHR Database



领域 Domain	映射逻辑 Mapping logic	映射率 (记录层级) Mapping rate (records level)
<b>诊断</b> (住院+门诊) <b>Diagnosis</b> (inpatient + outpatient)	<ul style="list-style-type: none"> <li>原始诊断名称首先通过中文名称匹配映射至 ICD10CN 编码。</li> <li>对于包含多个诊断的原始诊断，在适用的情况下将其映射至多个类型域。</li> <li>所有已映射的诊断名称最终将统一标准化为 OMOP 的 condition 类型域。</li> <li>Diagnosis names mapped to ICD10CN and OMOP concepts.</li> <li>Multiple diagnoses mapped to multiple concepts as needed.</li> <li>Standardized to OMOP concepts.</li> </ul>	89%
<b>用药</b> (住院+门诊) <b>Drug</b> (inpatient + outpatient)	<ul style="list-style-type: none"> <li>原始药品名称通过中文名称匹配映射至 ATC 编码</li> <li>所有已映射的药物名称最终将统一标准化为 OMOP 的 drug ingredient 类型域。</li> <li>Drug names mapped to ATC and OMOP concepts.</li> <li>Standardized to OMOP drug ingredients.</li> </ul>	85%
<b>手术</b> <b>Procedure</b>	<ul style="list-style-type: none"> <li>原始手术名称通过中文名称匹配映射至 ICD9ProcCN 编码；对于未能匹配的名称，在可行的情况下进一步映射至 OMOP 概念。</li> <li>所有已映射的手术名称最终将统一标准化为 OMOP 的 procedure 类型域。</li> <li>Operation names mapped to ICD9ProcCN via Chinese name matching; remaining names mapped to OMOP where possible.</li> <li>Standardized to OMOP procedure concepts.</li> </ul>	89%
<b>实验室检测</b> <b>Lab tests</b>	<ul style="list-style-type: none"> <li>原始实验室检测名称首先通过中文名称匹配映射至 LOINC CN 编码；对于未能匹配的名称，在可行的情况下进一步映射至 OMOP 概念。</li> <li>所有已映射的实验室检测名称最终将统一标准化为 OMOP 的 measurement 类型域。</li> <li>Lab names mapped to LOINC CN via Chinese name matching; remaining names mapped to OMOP where possible.</li> <li>Standardized to OMOP measurement concepts.</li> </ul>	87%
<b>实验室检测单位</b> <b>Lab test units</b>	<ul style="list-style-type: none"> <li>原始实验室检测单位将直接映射至 OMOP 的 unit 类型域。</li> <li>Raw lab test units were mapped directly to OMOP unit concepts.</li> </ul>	98%

# 天津电子健康档案数据库成功完成OMOP-CDM映射和质量控制

## Successful OMOP-CDM Mapping and Quality Assurance in Tianjin rEHR Database

领域 Domain	映射率 (记录层级) Mapping rate (records level)
诊断 (住院+门诊) Diagnosis (inpatient + outpatient)	89%
用药 (住院+门诊) Drug (inpatient + outpatient)	85%
手术 Procedure	89%
实验室检测 Lab tests	87%
实验室检测单位 Lab test units	98%



天津电子健康档案数据库 (rEHR) **超过85%的临床数据**已成功映射至 OMOP-CDM 结构  
**Over 85% of clinical records** in the Tianjin rEHR database were successfully mapped to the OMOP-CDM structure.



转换后的数据库通过**126项质量检查**, 涵盖以下领域的验证:

The transformed database **passed all 126 quality checks**, including validation of the following domains:

- 患者信息 Patient information
- 诊断 Diagnoses
- 药物 Medications
- 手术 Procedures
- 实验室检测 Laboratory measurements



## 前列腺癌患者的特征分析

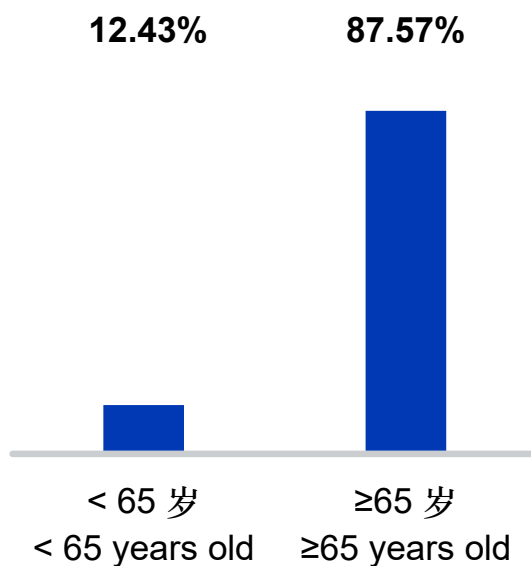
Characteristics of prostate cancer patients



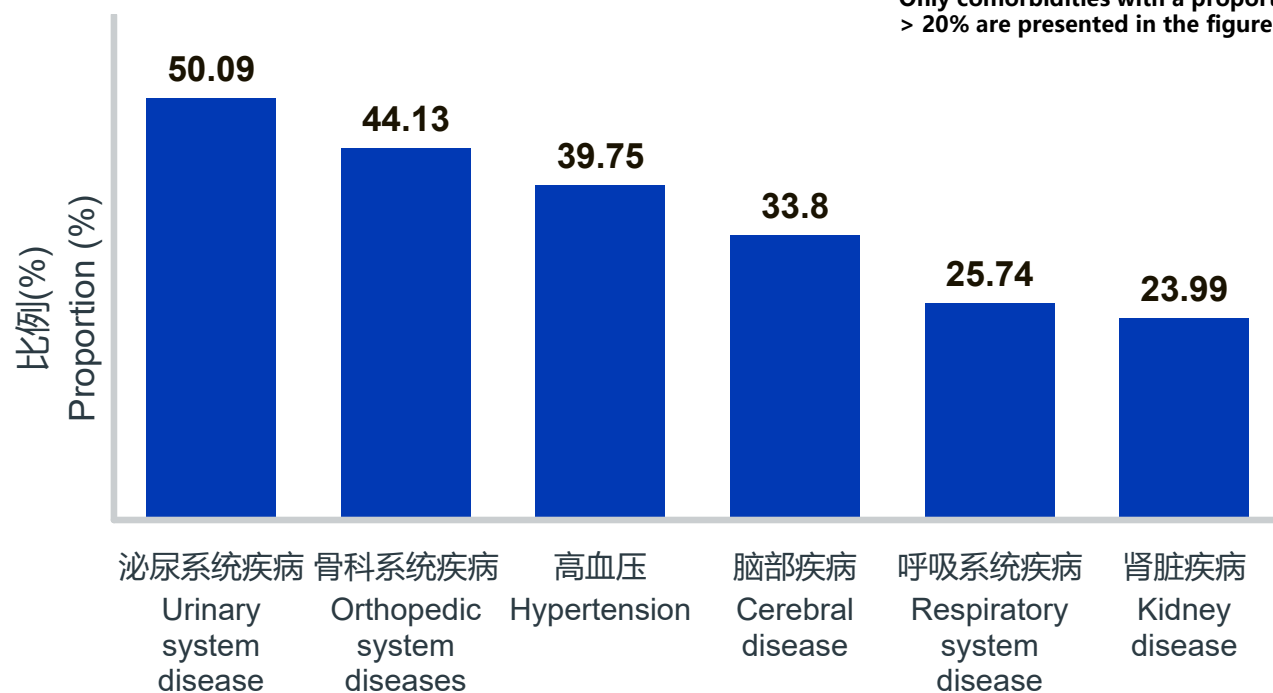
4<sup>th</sup> BOAO INTERNATIONAL CONFERENCE ON  
REAL WORLD STUDIES  
OF MEDICAL PRODUCTS

第四届博鳌国际药械真实世界研究大会

诊断时的年龄组  
Age group at diagnosis



共病\*  
Comorbidity\*



\* 图中仅展示比例 > 20% 的共病情况  
Only comorbidities with a proportion  
> 20% are presented in the figure.





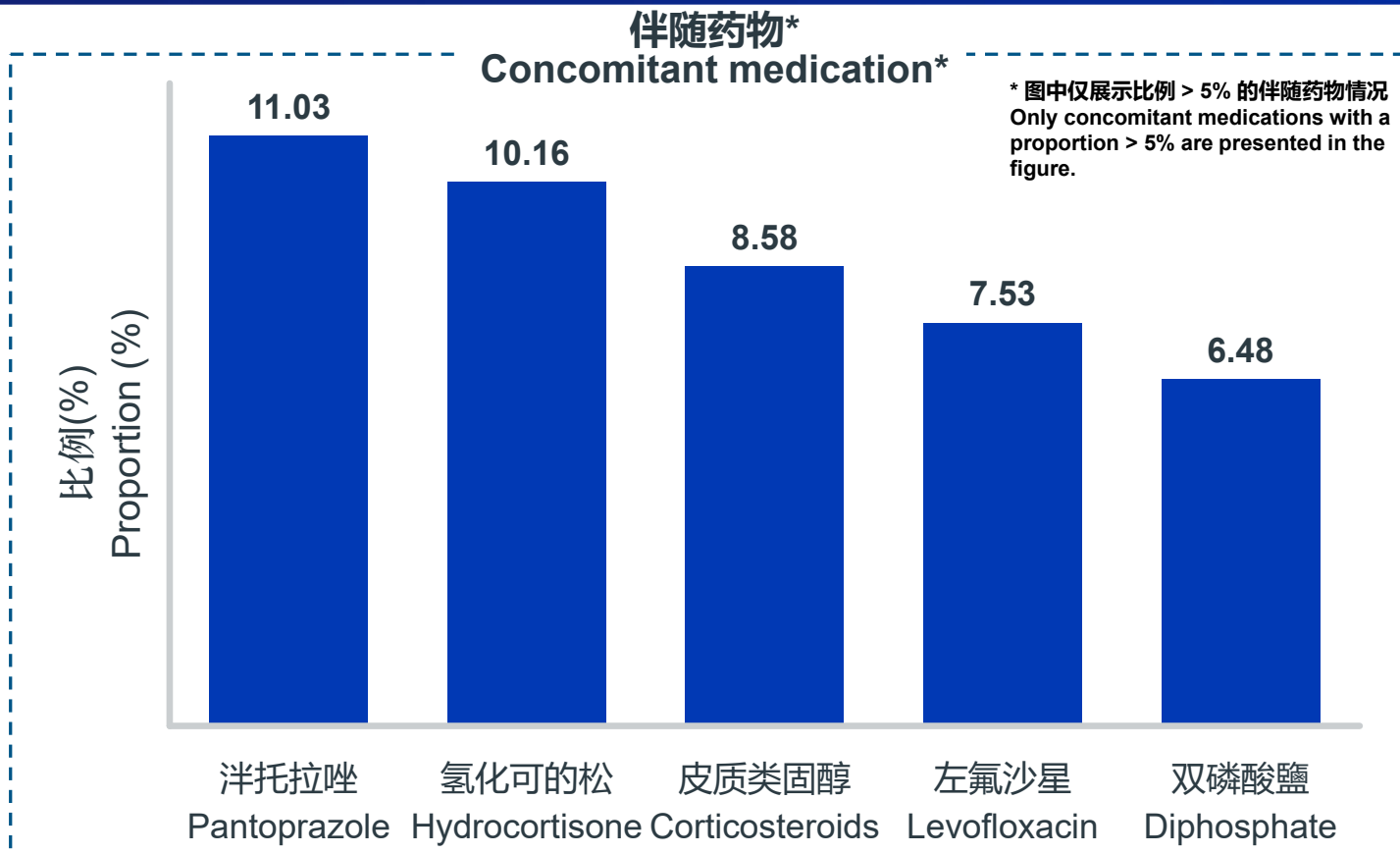
## 前列腺癌患者的特征分析

Characteristics of prostate cancer patients



4<sup>th</sup> BOAO INTERNATIONAL CONFERENCE ON  
REAL WORLD STUDIES  
OF MEDICAL PRODUCTS

第四届博鳌国际药械真实世界研究大会



## 天津电子健康档案数据库的前列腺癌队列已成功转换为 OMOP-CDM

The prostate cancer cohort of the Tianjin rEHR database has been successfully converted to the OMOP-CDM



本研究表明中国的区域电子健康档案数据可在**信息损失最小**的情况下实现标准化

This study demonstrated that standardizing regional EHR data in China is feasible with **minimal information loss**.



此次数据转换凸显了 **OMOP-CDM** 在中国的其他区域电子健康档案数据库或独立电子病历数据库中**应用的巨大潜力**

This success highlights the **potential to scale the OMOP-CDM approach** to other regional or standalone EMR systems across China.



这一标准化工作有助于**提升数据互操作性**，符合国家健康数据整合的目标，并进一步**强化了在本研究中应用该框架的价值**

This standardization effort supports **enhanced data interoperability**, aligns with national health data integration goals, and **reinforces the value of applying this framework in the current study**.



4<sup>th</sup>

# 第四届博鳌国际药械真实世界研究大会

BOAO INTERNATIONAL CONFERENCE ON **REAL WORLD STUDIES** OF MEDICAL PRODUCTS

感谢聆听!  
Thank you!

姓名：许哲 单位：拜耳医药保健有限公司

Presenter: Zhe Xu Affiliation: Bayer Healthcare Company Ltd

